

Prediction of PIK3CA mutation with gene expression

Jun Kang *

* Corresponding author: jkang.alien@gmail.com

Introduction

- Brevity
- Logic and clarity
- Clean typing
- The problem

– PIK3CA mutation in selecting drug

Targeted therapy becomes standard treatment in many cancer patients. To treat with the targeted therapy, the patients should be confirmed the indicated genomic alterations in their cancer. Many direct test for the genomic alteration has been developed and proved for their clinical utility to discriminate which patients will be response the targeted therapy.

The machine learning approach has been actively researched recently to detect the genomic alterations. The machine learning can build a prediction model from a large number of predictors such as radiomic features [1], pathology image or gene expression data [2]. The machine learning approach has limitations to apply to clinical practice. However the machine learning prediction can be a second best when the direct test fails.

PIK3CA encodes the p110 α catalytic subunit of phosphatidylinositol 3'-kinase (PI3K). PI3K is a protein kinase which phosphorylates phosphatidylinositol 4,5-bisphosphate (PIP₂) to make phosphatidylinositol 3,4,5-triphosphate (PIP₃). Phosphatase and tensin homolog (PTEN) changes PIP₂ to PIP₃ in contrast PI3K. PIP₃ is a second messenger to activate protein kinase B (AKT) which is a serine/threonine-specific protein kinase. AKT inhibits apoptosis and promotes cell proliferation.

Breast cancer with PIK3CA mutation has been approved to use PIK3CA inhibitor in hormone receptor positive HER2 negative subtype. [3] The PIK3CA mutation is second most driver mutation after TP53. The PIK3CA mutation is most frequently found in endometrial carcinoma (45%), and followed by breast invasive carcinoma (24%), cervical squamous cell carcinoma and endocervical adenocarcinoma (20%) and colon adenocarcinoma (16%).

Prediction RAS pathway activation with gene expression data was done in previous study. [4] They trained pancreatic The Cancer Genome Atlas (TCGA) data with a supervised elastic net penalized logistic regression classifier with stochastic gradient descent. The performance of their model was 84% with an area under the receiver operating characteristic (AUROC) curve and 63% with an area under the precision recall (AUPR) curve. We applied their modeling methods in prediction PIK3CA mutation.

Materials and Methods

Dataset

We used TCGA pancancer dataset. PIK3CA mutation data was get using cgdscr rpackage.[5] Gene expression data was get from GDAC firehose using RTCGAToolbox R package. [6]

10845 cases were available both PIK3CA mutation and mRNA expression data. 5128 out of 20502 genes were included in the modeling process after filtering with median absolute deviation as described at modelinf process method. 33 cancer type dummy variables were included in predictor variables.

The target variable was PIK3CA mutation status. The status of PIK3CA was considered as positive when the case has following PIK3CA variants which is the target variables of the thescreen PIK3CA RGQ PCR Kit; C420R, E542K, E545A, E545D, E545G, E545K, Q546E, Q546R, H1047L, H1047R, H1047Y. The thescreen PIK3CA RGQ PCR Kit was approved by U.S. Food and Drug adminidtration for PIK3CA inhibitor. We splited the three quarters of dataset for the trainset and one quarter for testset.

Modeling process

To narrow down potential predictors, Genes with a large the median absolute deviation (more than third-quartiles) were selected. Yeo-Johnson transformation was done to correct skewness. Centering and scaling were done. All preprocessing was done using recipe r package. [7] Penalized logistic regression was applied to prediction modeling. 10-fold cross-validation with targe variable stratification was done over the hyperparameter grid: $\lambda \{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^0\}$, $\alpha \{0.0, 0.25, 0.5, 0.75\}$. Lambda is penalty scaling parameter and alpha is mixing parameter of penalty function $((1 - \alpha)/2\|\beta\|_2^2 + \alpha\|\beta\|_1)$. [8] Model performance was eevaluated with the area under receiver operating characteristic (ROC) and area under the ROC (AUROC) and precision recall (PR) and area under curve (AUPR). The modelin process was done with tidymodels rpackage. [9]

Results

Dataset summary

Prevalence of PIK3CA was 0.11 in all cases. The PIK3CA prevalence of each cancer type was vary. The median prevalence of PIK3CA of each cancer types was 0.03 (range 0-0.33) (Figure 1).

Selecting model and performance estimation

In 10-fold cross-validation, the model with $\lambda = 0.01$ and $\alpha = 1.0$ (Ridge regression) were best performance in terms of AUROC. The final model was trained with the selected hyperparameters with all trainset.

The trainset AUROC was 0.93 and the testset AUROC was 0.84. The AUPR of trainset was 0.66 and the testset AUPR was 0.39. (Figure 1A)

Performance of each cancer type

Because tht prevalence of PIK3CA mutation is vary across the cancer type, the performance of each cancer type was investigated. The AUROC and AUPR were

positively correlated between train set and test set in cancer type subanalysis. (Figure 1B) The AUPR was high in cancer type with high PIK3CA mutation rate such as colon, breast, Uterus cancer types. The AUROC did not correlated with PIK3CA mutation rate of each cancer types. (Figure 1C)

Important predictors

Figure 2 shows top 30 important predictors. The coefficient is the parameter of the predictor which represent the effect of the predictor on prediction. IGF1R mRNA expression was the strongest negative predictor and PTEN was the strongest positive predictor. Both genes and PIK3CA are key players in tyrosin kanase pathway. The cancer type was important predictors. Some cancer types including uterine carcinosarcoma (UCS), bladder urothelial carcinoma (BLCA), pancreatic adenocarcinoma (PAAD), lymphoid neoplasm diffuse large B-cell lymphoma (DLBC) are strongest predictors.

Discussion

Our model showed good performance to predict PIK3CA mutation of various cancer types.

- Main message answers the question and main supporting evidence
- Critical assessment opinions on any shortcomings in study design

In this study, the split method was done for testset instead of indepedent testset. Because there is no standard analysis method for the RNA-seq gene expression quantification, it is difficult to find gene expression data same as TGCA data. Developing standard analysis method for gene expression data is neccessory to apply this prediction model.

In this study, the prevalence of PIK3CA mutation was 11%. The low prevalence of PIK3CA results in imbalaced dataset which makes the prediction difficult. The AUROC is optimistic in terms of performance. The AUPR is more informative than AUROC on imbalaced datasets. [10]

- limitations in methods
 - Case imbalance
 - flaws in analysis
 - validity of assumption
- Comparison with other studies where inconsistencies are discussed
- Evaluate the results - not the authors

Our PIK3CA prediction model performed better than RAS activation prediction model of previous study in terms of both AUROC (0.84 vs 0.75) and AUPR (0.39 vs 0.24) on testset prediction. Our testset is corresponding to the samples initially filtered from training. The target variable of our study is more specific than the previous study. The specific important mutations can effect stronger downsteam gene expression than the broad events pathway activation. The previous study might be more difficult prediction problem.

Our model includes cancer type predictor and they are stronger than gene expression data. The varing prevalence of PIC3CA mutation across cancer type might reason of the strong cancer type predictor. If the cancer type was wrong or can not be determined, our model performance can be poor.

- Conclusions comments on possible biological or clinical implications and suggestions for further research. 122
- Proof of concept study 123
- Reproducibility of gene expression prediction model 124

Figure legends 125

- Figure 1. Prevalance rate of PIK3CA across cancer types The abbreviations of cancer types are explained in S1 appendix. 126
- Figure 2. Summary of modeling results 127
- (A) Left: receiver operating characteristic (ROC) curve right: precision recall (PR) curve of trainset and testset. The horizontal green line is the PIK3CA mutation rate (0.11) (B) Correlation between trainset and testset of area under receiver operating characteristic curve (AUROC) and area under precision recall curve (AUPR) among cancer type. The abbreviations are explained in S1 appendix. (C) Correlation between the PIK3CA mutation rate of area under receiver operating characteristic curve (AUROC) and area under precision recall curve (AUPR). 128
- Figure 3. Coefficients of model 129
- (A) Top 30 high coefficients of mRNA. (B) Coefficients of cancer types. The abbreviations of cancer types are explained in S1 appendix. 130

Supporting information 131

- S1 Appendix. 132
- S2 Figure. 133
- S1 Table. 134
- S2 Table. 135
- S3 Table. 136

References 137

1. Dercle L, Fronheiser M, Lu L, Du S, Hayes W, Leung DK, et al. Identification of NonSmall Cell Lung Cancer Sensitive to Systemic Cancer Therapies Using Radiomics. Clin Cancer Res. American Association for Cancer Research; 2020;26: 2151–2162. doi:10.1158/1078-0432.CCR-19-2942 138
2. Way GP, Sanchez-Vega F, La K, Armenia J, Chatila WK, Luna A, et al. Machine Learning Detects Pan-cancer Ras Pathway Activation in The Cancer Genome Atlas. Cell Reports. Elsevier; 2018;23: 172–180.e3. doi:10.1016/j.celrep.2018.03.046 139
3. André F, Ciruelos E, Rubovszky G, Campone M, Loibl S, Rugo HS, et al. Alpelisib for PIK3CA-Mutated, Hormone ReceptorPositive Advanced Breast Cancer. New England Journal of Medicine. Massachusetts Medical Society; 2019;380: 1929–1940. doi:10.1056/NEJMoa1813904 140
4. Way GP, Sanchez-Vega F, La K, Armenia J, Chatila WK, Luna A, et al. Machine Learning Detects Pan-cancer Ras Pathway Activation in The Cancer Genome Atlas. Cell Reports. Elsevier; 2018;23: 172–180.e3. doi:10.1016/j.celrep.2018.03.046 141
5. Jacobsen A, Luna A. Cgdsr: R-based API for accessing the MSKCC cancer genomics data server (CGDS). 2019. 142

6. Samur MK. RTCGAToolbox: A new tool for exporting TCGA Firehose data.	163
PLoS One. 2014;9(9):e106397.	164
7. Kuhn M, Wickham H. Recipes: Preprocessing tools to create design matrices.	165
2020.	166
8. Friedman J, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear	167
Models via Coordinate Descent. J Stat Softw. 2010;33: 1–22.	168
9. Kuhn M, Wickham H. Tidymodels: Easily install and load the 'tidymodels'	169
packages. 2020.	170
10. Saito T, Rehmsmeier M. The Precision-Recall Plot Is More Informative than the	171
ROC Plot When Evaluating Binary Classifiers on Imbalanced Datasets. PLOS ONE.	172
Public Library of Science; 10: e0118432. doi:10.1371/journal.pone.0118432	173